

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号
特開2001-229026
(P2001-229026A)

(43) 公開日 平成13年8月24日 (2001.8.24)

(51) Int.Cl. ⁷	識別記号	F I	メモコード* (参考)
G 0 6 F 9/44	5 8 0	G 0 6 F 9/44	5 8 0 A
G 0 6 N 1/00		G 0 6 N 1/00	

審査請求 有 請求項の数15 O L (全 10 頁)

(21) 出願番号 特願2000-104760 (P2000-104760)

(22) 出願日 平成12年4月6日 (2000.4.6)

(31) 優先権主張番号 特願平11-350834

(32) 優先日 平成11年12月9日 (1999.12.9)

(33) 優先権主張国 日本 (J P)

(71) 出願人 000004237

日本電気株式会社
東京都港区芝五丁目7番1号

(72) 発明者 馬見塚 拓

東京都港区芝五丁目7番1号 日本電気株
式会社内

(72) 発明者 安倍 直樹

東京都港区芝五丁目7番1号 日本電気株
式会社内

(74) 代理人 100108578

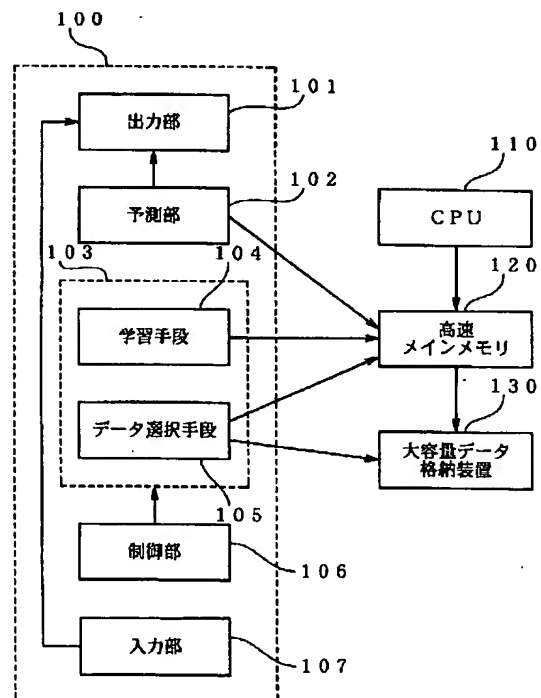
弁理士 高橋 詔男 (外3名)

(54) 【発明の名称】 知識発見方式

(57) 【要約】

【課題】 データベースから情報量の大きいデータのみ選択的にサンプリングし、効率的に知識発見を行なう汎用的な知識発見方式を提供する。

【解決手段】 学習手段104は、高速メインメモリ120に格納されているデータからサンプリングによって作成した複数の部分サンプルを、入力部107を介して入力した下位学習アルゴリズムに学習させ、複数の仮説を得る。データ選択手段105は、この仮説を用いて、大容量データ格納装置130より読み出された候補データ各々の情報量を推定し、情報量の大きいデータのみを高速メインメモリ120に追加格納する。制御部106は、上記の処理を所定の回数繰り返し、得られた最終仮説を格納する。そして、予測部102は入力部107へ入力されたラベル未知のデータに対し最終仮説によりラベル値を予測し、出力部101はこの予測値を出力する。



【特許請求の範囲】

【請求項1】 大容量データ格納装置に格納されたデータベースからサンプリングしたデータを計算機のメインメモリに読み込み高次知識を抽出するデータマイニングの知識発見方式において、

学習アルゴリズムおよび該学習アルゴリズムに学習させるデータを入力する入力手段と、

前記メインメモリに格納されているデータからサンプリングして作成した複数の部分データ集合を前記学習アルゴリズムへ訓練データとして入力して学習させ、複数の仮説を得る学習手段と、

該学習手段により得られた複数の仮説を用いて前記大容量データ格納装置から読み出された複数のサンプル候補点に対する関数値の予測を行い、求めた予測値に基づき候補点の情報量を推定し、情報量の大きい候補点を1つあるいは複数選択し、前記メインメモリに追加して格納するデータ選択手段と、

前記学習手段と前記データ選択手段による選択的なデータ格納および知識発見の処理を予め定めた停止条件が満たされるまで繰り返し、その結果得られた複数の仮説を最終仮説として前記メインメモリに格納させる制御手段と、

前記複数の仮説間の平均または重み付き平均をもって前記入力部に入力したラベル未知のデータに対してラベル値を予測する予測手段と、

を有することを特徴とする知識発見方式。

【請求項2】 前記データ選択手段は、前記学習手段により得られた複数の仮説を用いて前記サンプル候補点に対する関数値の予測を行い、求めた予測値の分散値により前記候補点の情報量を推定し、分散値の大きい候補点を1つあるいは複数選択し、前記メインメモリに追加して格納することを特徴とする請求項1に記載の知識発見方式。

【請求項3】 前記データ選択手段は、前記学習手段により得られた各仮説の前記訓練データに対する予測誤差の関数として各仮説の重みを計算し、データ選択の際に、多値予測の場合、各候補点に対して予測値を求める前記仮説の重みの総和を算出し、最も大きい重みの総和と次に大きい重みの総和の差分であるマージンを求め、このマージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記仮説による予測値の重み付き分散値の大きい候補点を1つまたは複数選択して前記メインメモリに追加して格納し、前記予測手段は、前記仮説の重みによる重み付き平均をもって最終仮説による予測を行うことを特徴とする請求項1に記載の知識発見方式。

【請求項4】 前記データ選択手段は、前記大容量データ格納装置から読み出されたデータに含まれるデータ候補点の真のラベルを利用して、多値予測の場合、前記仮説の予測誤差を用いて誤差マージンを算出し、該誤差マ

ージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記予測値の重み付き予測誤差の大きい候補点を1つないし複数選択して前記メインメモリに追加して格納することを特徴とする請求項1に記載の知識発見方式。

【請求項5】 前記データ選択手段は、前記候補点を選択し前記メインメモリに格納するとき、既に前記メインメモリに格納されているデータに対して、前記複数の仮説を用いて関数値の予測値を求め、前記予測値の分散の小さい候補点を1つまたは複数選択して前記メインメモリの格納データから削除することを特徴とする請求項1または2に記載の知識発見方式。

【請求項6】 前記データ選択手段は、前記候補点を選択し前記メインメモリに格納するとき、既に前記メインメモリに格納されているデータに対して、前記複数の仮説を用いて関数値の予測値を求め、多値予測の場合、前記予測値のマージンの大きい候補点を1つまたは複数選択し、実数値予測の場合、前記予測値の重み付き分散値の小さい候補点を1つまたは複数選択し、前記メインメモリの格納データから削除することを特徴とする請求項1または3に記載の知識発見方式。

【請求項7】 前記データ選択手段は、前記候補点を選択し前記メインメモリに格納するとき、既に前記メインメモリに格納されているデータに対して、前記複数の仮説を用いて関数値の予測値を求め、多値予測の場合、前記予測値の誤差マージンの大きい候補点を1つまたは複数選択し、実数値予測の場合、前記予測値の重み付き予測誤差の小さい候補点を1つまたは複数選択し、前記メインメモリの格納データから削除することを特徴とする請求項1または4に記載の知識発見方式。

【請求項8】 大容量データ格納装置に格納されたデータベースからサンプリングしたデータを計算機のメインメモリに読み込み高次知識を抽出するデータマイニングの知識発見方式において、

学習アルゴリズムおよび該学習アルゴリズムに学習させるデータを入力する入力手段と、

前記メインメモリに格納されたデータを訓練データとして入力して学習し、学習した仮説を大容量データ格納装置もしくはメインメモリに格納する学習手段と、

該学習手段により得られた過去の複数の仮説を用いて、前記大容量データ格納装置から読み出された複数のサンプル候補点に対する関数値の予測を行い、求めた予測値に基づき候補点の情報量を推定し、情報量の大きい候補点を複数選択し、前記メインメモリに格納するデータ選択手段と、

前記学習手段と前記データ選択手段におけるデータ格納および知識発見の処理を予め定めた停止条件が満たされるまで繰り返し、その結果得られた複数の仮説を最終仮説として前記メインメモリもしくは大容量データ格納装置に格納させる制御手段と、

前記複数の仮説間の平均または重み付き平均をもって前記入力部に入力したラベル未知のデータに対してラベル値を予測する予測手段と、

を有することを特徴とする知識発見方式。

【請求項 9】 前記データ選択手段は、前記学習手段により得られた複数の仮説を用いて前記サンプル候補点に対する関数値の予測を行い、求めた予測値の分散値により前記候補点の情報量を推定し、分散値の大きい候補点を 1 つあるいは複数選択し、前記メインメモリに格納することを特徴とする請求項 8 に記載の知識発見方式。

【請求項 10】 前記データ選択手段は、前記学習手段により得られた各仮説の前記訓練データに対する予測誤差の関数として各仮説の重みを計算し、データ選択の際に、多値予測の場合、各候補点に対して予測値を求める前記仮説の重みの総和を算出し、最も大きい重みの総和と次に大きい重みの総和の差分であるマージンを求め、このマージンの小さい候補点を 1 つまたは複数選択し、あるいは、実数値予測の場合、前記仮説による予測値の重み付き分散値の大きい候補点を 1 つまたは複数選択して前記メインメモリに格納し、

前記予測手段は、前記仮説の重みによる重み付き平均をもって最終仮説による予測を行うことを特徴とする請求項 8 に記載の知識発見方式。

【請求項 11】 前記データ選択手段は、前記大容量データ格納装置から読み出されたデータに含まれるデータ候補点の真のラベルを利用して、多値予測の場合、前記仮説の予測誤差を用いて誤差マージンを算出し、該誤差マージンの小さい候補点を 1 つまたは複数選択し、あるいは、実数値予測の場合、前記予測値の重み付き予測誤差の大きい候補点を 1 つないし複数選択して前記メインメモリに格納することを特徴とする請求項 8 に記載の知識発見方式。

【請求項 12】 大容量データ格納装置に格納されたデータベースからサンプリングしたデータを計算機のメインメモリに読み込み高次知識を抽出するデータマイニングの知識発見方式において、

学習アルゴリズムおよび該学習アルゴリズムに学習させるデータを入力する入力手段と、

前記メインメモリに格納されているデータからサンプリングして作成した複数の部分データ集合を前記学習アルゴリズムへ訓練データとして入力して学習させ、複数の仮説を得る学習手段と、

該学習手段により得られた過去の複数の仮説を用いて、前記大容量データ格納装置から読み出された複数のサンプル候補点に対する関数値の予測を行い、求めた予測値に基づき候補点の情報量を推定し、情報量の大きい候補点を複数選択し、前記メインメモリに格納するデータ選択手段と、

前記学習手段と前記データ選択手段におけるデータ格納および知識発見の処理を予め定めた停止条件が満たされ

るまで繰り返し、その結果得られた複数の仮説を最終仮説として前記メインメモリもしくは大容量データ格納装置に格納させる制御手段と、

前記複数の仮説間の平均または重み付き平均をもって前記入力部に入力したラベル未知のデータに対してラベル値を予測する予測手段と、

を有することを特徴とする知識発見方式。

【請求項 13】 前記データ選択手段は、前記学習手段により得られた複数の仮説を用いて前記サンプル候補点に対する関数値の予測を行い、求めた予測値の分散値により前記候補点の情報量を推定し、分散値の大きい候補点を 1 つあるいは複数選択し、前記メインメモリに格納することを特徴とする請求項 12 に記載の知識発見方式。

【請求項 14】 前記データ選択手段は、前記学習手段により得られた各仮説の前記訓練データに対する予測誤差の関数として各仮説の重みを計算し、データ選択の際に、多値予測の場合、各候補点に対して予測値を求める前記仮説の重みの総和を算出し、最も大きい重みの総和と次に大きい重みの総和の差分であるマージンを求め、このマージンの小さい候補点を 1 つまたは複数選択し、あるいは、実数値予測の場合、前記仮説による予測値の重み付き分散値の大きい候補点を 1 つまたは複数選択して前記メインメモリに格納し、前記予測手段は、前記仮説の重みによる重み付き平均をもって最終仮説による予測を行うことを特徴とする請求項 12 に記載の知識発見方式。

【請求項 15】 前記データ選択手段は、前記大容量データ格納装置から読み出されたデータに含まれるデータ候補点の真のラベルを利用して、多値予測の場合、前記仮説の予測誤差を用いて誤差マージンを算出し、該誤差マージンの小さい候補点を 1 つまたは複数選択し、あるいは、実数値予測の場合、前記予測値の重み付き予測誤差の大きい候補点を 1 つないし複数選択して前記メインメモリに格納することを特徴とする請求項 12 に記載の知識発見方式。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、大容量ディスクなどのデータ格納装置に格納されたデータベースから有益な知識を抽出するデータマイニングに用いて好適な知識発見方式に関する。

【0002】

【従来の技術】従来、データマイニングまたは知識発見方式は、ディスクに蓄えられた大量のデータに対してなるべく少回数アクセスし、読み込むことにより有効な知識発見を行なうことに主眼が置かれていた。例えば、結合ルールの抽出方式がその代表的なものであり、この方式については 1993 年発行の国際会議の予稿集「プロシーディングス オブ エーシーエム シングモッド

コンファレンス オンマネージメント オブ データ (Proceedings of ACM SIGMOD Conference on Management of data) 2007頁-216頁に記載のアグラワル (Agrawal)らによる論文「マイニング アソシエーション ルールズ ビトウィーン セッツ オブ アイテムズ イン ラージ データベース (Mining Association Rules Between Sets of Items in Large Databases)」に詳細に記載されている。

【0003】また、ディスクに蓄えられたデータの選択的サンプリングの方式としては、集団質問学習と呼ばれる方法が知られており、1992年発行の国際会議の予稿集「プロシーディングス オブ ザ フィフス アニュアル エーシーエム ワークショップ オン コンピュータショナル ラーニング セオリー (Proceedings of the fifth annual ACM workshop on computational learning theory)」287頁-294頁に記載のセング (Seung)らの論文「クエリー バイ コミッティ (Query by Committee)」に詳細が記載されている。この方法では、ランダム化された下位学習方式に複数回予測をさせ、その不一致度が大きい点のみを選択するというを行う。この方法に用いられる下位学習方式はランダム化された高性能な学習方式であることが前提になっている。

【0004】一方、性能の比較的低い学習方式の精度を増強するという文脈において、与えられたデータから繰り返し再サンプリングを行ない、そのデータを用いて学習させて得られた複数の仮説を統合することにより学習精度を向上させる一連の技術が近年注目されている。こうした技術の代表的な手法はバグギング並びにブースティングであり、これらの技術は実験的に高い性能を有することが確かめられている。バグギングの手法については、1994年発行のカリフォルニア大学バークレー校の技術報告書421に記載されたブライマン (Breiman)の論文「バグギング プレディクターズ (Bagging Predictors)」に記載されている。

【0005】ブースティングの手法については、1995年発行の国際会議予稿集「プロシーディングス オブ ザ セカンド ヨーロピアン コンファレンス オン コンピュータショナル ラーニング セオリー (Proceedings of the second european conference on computational learning theory)」23頁-37頁に記載のフロインド (Freund)とシャピレ (Shapire)の論文「ア デジジョン セオレティック ジェネライゼーション オブ オンライン ラーニング アンド アンアプリケーション トゥー ブースティング (A decision-theoretic generalization of on-line learning and an application to boosting)」に記載されている。上記の集団質問学習方式は、テキスト分類等の分野において、人手で分類しラベルづける文書を選択する問題等に適用されている。また、バグギング技術やブースティング技術は、受動学習における精度増強の目的で利用されてきて

いる。

【0006】

【発明が解決しようとする課題】ところで、上述したように、従来のデータマイニング手法では、効率的なデータマイニングのために、選択的サンプリング方式を用いる場合、下位学習方式はランダム化された高性能な学習方式でなければならないという問題がある。また、バグギング技術あるいはブースティング技術を用いて性能の比較的低い学習方式の精度を上げることはできるが、両方の技術の利点を合わせて用いることができる選択的サンプリング方式がないという問題がある。

【0007】この発明は、上記の点に鑑みてなされたもので、その目的は、大容量データ格納装置に格納された大量データからサンプル候補データを読み出し、比較的低い学習方式を用いて学習させ、その結果により情報量の大きいデータを選択してメインメモリに読み込むことにより、効率的に高精度のラベル予測規則を発見する汎用的な知識発見方式を提供することにある。

【0008】

【課題を解決するための手段】上記の目的を達成するために、請求項1に記載の発明は、大容量データ格納装置に格納されたデータベースからサンプリングしたデータを計算機のメインメモリに読み込み高次知識を抽出するデータマイニングの知識発見方式において、学習アルゴリズムおよび該学習アルゴリズムに学習させるデータを入力する入力手段と、前記メインメモリに格納されているデータからサンプリングして作成した複数の部分データ集合を前記学習アルゴリズムへ訓練データとして入力して学習させ、複数の仮説を得る学習手段と、該学習手段により得られた複数の仮説を用いて前記大容量データ格納装置から読み出された複数のサンプル候補点に対する関数値の予測を行い、求めた予測値に基づき候補点の情報量を推定し、情報量の大きい候補点を1つあるいは複数選択し、前記メインメモリに追加して格納するデータ選択手段と、前記学習手段と前記データ選択手段による選択的なデータ格納および知識発見の処理を予め定めた停止条件が満たされるまで繰り返し、その結果得られた複数の仮説を最終仮説として前記メインメモリに格納させる制御手段と、前記複数の仮説間の平均または重み付き平均をもって前記入力部に入力したラベル未知のデータに対してラベル値を予測する予測手段とを有することを特徴とする。

【0009】請求項2に記載の発明は、請求項1に記載の知識発見方式において、前記データ選択手段は、前記学習手段により得られた複数の仮説を用いて前記サンプル候補点に対する関数値の予測を行い、求めた予測値の分散値により前記候補点の情報量を推定し、分散値の大きい候補点を1つあるいは複数選択し、前記メインメモリに追加して格納することを特徴とする。請求項3に記載の発明は、請求項1に記載の知識発見方式において、

前記データ選択手段は、前記学習手段により得られた各仮説の前記訓練データに対する予測誤差の関数として各仮説の重みを計算し、データ選択の際に、多値予測の場合、各候補点に対して予測値を求める前記仮説の重みの総和を算出し、最も大きい重みの総和と次に大きい重みの総和の差分であるマージンを求め、このマージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記仮説による予測値の重み付き分散値の大きい候補点を1つまたは複数選択して前記メインメモリに追加して格納し、前記予測手段は、前記仮説の重みによる重み付き平均をもって最終仮説による予測を行うことを特徴とする。

【0010】請求項4に記載の発明は、請求項1に記載の知識発見方式において、前記データ選択手段は、前記大容量データ格納装置から読み出されたデータに含まれるデータ候補点の真のラベルを利用して、多値予測の場合、前記仮説の予測誤差を用いて誤差マージンを算出し、該誤差マージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記予測値の重み付き予測誤差の大きい候補点を1つないし複数選択して前記メインメモリに追加して格納することを特徴とする。請求項5に記載の発明は、請求項1または2に記載の知識発見方式において、前記データ選択手段は、前記候補点を選択し前記メインメモリに格納するとき、既に前記メインメモリに格納されているデータに対して、前記複数の仮説を用いて関数値の予測値を求め、前記予測値の分散の小さい候補点を1つまたは複数選択して前記メインメモリの格納データから削除することを特徴とする。

【0011】請求項6に記載の発明は、請求項1または3に記載の知識発見方式において、前記データ選択手段は、前記候補点を選択し前記メインメモリに格納するとき、既に前記メインメモリに格納されているデータに対して、前記複数の仮説を用いて関数値の予測値を求め、多値予測の場合、前記予測値のマージンの大きい候補点を1つまたは複数選択し、実数値予測の場合、前記予測値の重み付き分散値の小さい候補点を1つまたは複数選択し、前記メインメモリの格納データから削除することを特徴とする。請求項7に記載の発明は、前記データ選択手段は、請求項1または4に記載の知識発見方式において、前記候補点を選択し前記メインメモリに格納するとき、既に前記メインメモリに格納されているデータに対して、前記複数の仮説を用いて関数値の予測値を求め、多値予測の場合、前記予測値の誤差マージンの大きい候補点を1つまたは複数選択し、実数値予測の場合、前記予測値の重み付き予測誤差の小さい候補点を1つまたは複数選択し、前記メインメモリの格納データから削除することを特徴とする。

【0012】請求項8に記載の発明は、大容量データ格納装置に格納されたデータベースからサンプリングした

データを計算機のメインメモリに読み込み高次知識を抽出するデータマイニングの知識発見方式において、学習アルゴリズムおよび該学習アルゴリズムに学習させるデータを入力する入力手段と、前記メインメモリに格納されたデータを訓練データとして入力して学習し、学習した仮説を大容量データ格納装置もしくはメインメモリに格納する学習手段と、該学習手段により得られた過去の複数の仮説を用いて、前記大容量データ格納装置から読み出された複数のサンプル候補点に対する関数値の予測を行い、求めた予測値に基づき候補点の情報を推定し、情報の大きい候補点を複数選択し、前記メインメモリに格納するデータ選択手段と、前記学習手段と前記データ選択手段におけるデータ格納および知識発見の処理を予め定めた停止条件が満たされるまで繰り返し、その結果得られた複数の仮説を最終仮説として前記メインメモリもしくは大容量データ格納装置に格納させる制御手段と、前記複数の仮説間の平均または重み付き平均をもって前記入力部に入力したラベル未知のデータに対してラベル値を予測する予測手段と、を有することを特徴とする。

【0013】請求項9に記載の発明は、請求項8に記載の知識発見方式において、前記データ選択手段は、前記学習手段により得られた複数の仮説を用いて前記サンプル候補点に対する関数値の予測を行い、求めた予測値の分散値により前記候補点の情報を推定し、分散値の大きい候補点を1つあるいは複数選択し、前記メインメモリに格納することを特徴とする。請求項10に記載の発明は、請求項8に記載の知識発見方式において、前記データ選択手段は、前記学習手段により得られた各仮説の前記訓練データに対する予測誤差の関数として各仮説の重みを計算し、データ選択の際に、多値予測の場合、各候補点に対して予測値を求める前記仮説の重みの総和を算出し、最も大きい重みの総和と次に大きい重みの総和の差分であるマージンを求め、このマージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記仮説による予測値の重み付き分散値の大きい候補点を1つまたは複数選択して前記メインメモリに格納し、前記予測手段は、前記仮説の重みによる重み付き平均をもって最終仮説による予測を行うことを特徴とする。

【0014】請求項11に記載の発明は、請求項8に記載の知識発見方式において、前記データ選択手段は、前記大容量データ格納装置から読み出されたデータに含まれるデータ候補点の真のラベルを利用して、多値予測の場合、前記仮説の予測誤差を用いて誤差マージンを算出し、該誤差マージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記予測値の重み付き予測誤差の大きい候補点を1つないし複数選択して前記メインメモリに格納することを特徴とする。

【0015】請求項12に記載の発明は、大容量データ

格納装置に格納されたデータベースからサンプリングしたデータを計算機のメインメモリに読み込み高次知識を抽出するデータマイニングの知識発見方式において、学習アルゴリズムおよび該学習アルゴリズムに学習させるデータを入力する入力手段と、前記メインメモリに格納されているデータからサンプリングして作成した複数の部分データ集合を前記学習アルゴリズムへ訓練データとして入力して学習させ、複数の仮説を得る学習手段と、該学習手段により得られた過去の複数の仮説を用いて、前記大容量データ格納装置から読み出された複数のサンプル候補点に対する関数値の予測を行い、求めた予測値に基づき候補点の情報量を推定し、情報量の大きい候補点を複数選択し、前記メインメモリに格納するデータ選択手段と、前記学習手段と前記データ選択手段におけるデータ格納および知識発見の処理を予め定めた停止条件が満たされるまで繰り返し、その結果得られた複数の仮説を最終仮説として前記メインメモリもしくは大容量データ格納装置に格納させる制御手段と、前記複数の仮説間の平均または重み付き平均をもって前記入力部に入力したラベル未知のデータに対してラベル値を予測する予測手段とを有することを特徴とする。

【0016】請求項13に記載の発明は、請求項12に記載の知識発見方式において、前記データ選択手段は、前記学習手段により得られた複数の仮説を用いて前記サンプル候補点に対する関数値の予測を行い、求めた予測値の分散値により前記候補点の情報量を推定し、分散値の大きい候補点を1つあるいは複数選択し、前記メインメモリに格納することを特徴とする。請求項14に記載の発明は、請求項12に記載の知識発見方式において、前記データ選択手段は、前記学習手段により得られた各仮説の前記訓練データに対する予測誤差の関数として各仮説の重みを計算し、データ選択の際に、多値予測の場合

$$S = \langle (x_1, y_1), \dots, (x_m, y_m) \rangle$$

ここで、 m はデータ数であり、 x_i はデータ中のあるレコードの予め指定された属性の組みの属性値の組み、 y_i は同一レコードの予め指定されたラベル属性の属性値である。

【0018】次に、上位アルゴリズムについて説明する。ここで、大容量データ格納装置から一部データを選択的にメモリに格納し、リサンプリング、および学習をするという手順を繰り返す回数（ステージ数）を N 、各ステージにおけるリサンプリングの回数を T 、メインメモリに格納すべきデータ点を選ぶ際の候補点の数を R 、その中から実際にメインメモリに格納する点の数を Q とする。上位アルゴリズムは、以下の手順1から手順3を N 回繰り返す。もしくは、手順1の代わりに手順1'、また、手順3の代わりに手順3'を使用する。

【0019】（手順1）メインメモリに格納されたデータからリサンプリングにより得られた複数のデータ集合 S_1, \dots, S_T に対して、下位学習アルゴリズムを走らせて

合、各候補点に対して予測値を求める前記仮説の重みの総和を算出し、最も大きい重みの総和と次に大きい重みの総和の差分であるマージンを求め、このマージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記仮説による予測値の重み付き分散値の大きい候補点を1つまたは複数選択して前記メインメモリに格納し、前記予測手段は、前記仮説の重みによる重み付き平均をもって最終仮説による予測を行うことを特徴とする。請求項15に記載の発明は、請求項12に記載の知識発見方式において、前記データ選択手段は、前記大容量データ格納装置から読み出されたデータに含まれるデータ候補点の真のラベルを利用して、多値予測の場合、前記仮説の予測誤差を用いて誤差マージンを算出し、該誤差マージンの小さい候補点を1つまたは複数選択し、あるいは、実数値予測の場合、前記予測値の重み付き予測誤差の大きい候補点を1つないし複数選択して前記メインメモリに格納することを特徴とする。

【0017】

【発明の実施の形態】以下、先ず、この発明の基本的考え方を説明する。この発明のアルゴリズムは、入力として与えられる下位学習アルゴリズムと、これを用いて、選択的サンプリングを行いながら知識発見を行なう上位アルゴリズムからなる。下位学習アルゴリズムの機能は、入力されたデータ S から学習を行い仮説を出力することと、データ S の1つのデータ点（属性値の組み） x に対して、学習により得られた仮説を使用して、そのラベル y の予測値を出力することである。ここで使用する学習アルゴリズムは、高度な学習性能を持つ学習アルゴリズムは必要なく、例えば、決定木を学習するアルゴリズムや階層型のニューラルネットワークの学習アルゴリズムなどを用いることができる。ラベル付きの学習データ S は次式で表される。

(1)

仮説 H_1, \dots, H_T を得る。ここで、データ $S = \langle (x_1, y_1), \dots, (x_m, y_m) \rangle$ からのリサンプリングとは、例えばデータ S から一様分布によりサンプリングする（即ち、データ S 中の各データを $1/m$ の確率で取り出す）ことを一定回数繰り返すことを言う。また、仮説とは、任意の属性値の組みに対して、そのラベルの予測値を規定するような表現形である。

（手順1'）メインメモリに格納されたデータ S に対して、下位アルゴリズム A を走らせて仮説 H を得る。

【0020】（手順2） R 個の候補サンプル点 x_1, \dots, x_R を大容量データ格納装置からメインメモリに読み込み、その中から情報量の大きい Q 個のサンプル点を選択しメインメモリに格納する。他の点についてはメインメモリから削除する。情報量の大きいサンプル点の選択においては、各候補点に対し、上記手順1で得られた T 個の仮説によりそのラベルの予測をさせて、その予測値の分散が大きい Q 個の点を選択する。

(手順3) 過去のデータに、手順2において得られたデータを加える。

(手順3') 過去のデータを、ステップ2において得られたデータに置き換える。

【0021】なお、上記のリサンプリングの手順を逐次的に変化する分布を用いて行なうことも可能である。例えば、精度増強法としてブースティング方式を用いる場合には、直前回の部分サンプルに対して得られた仮説の予測値が誤るようなデータの分布を逐次的に変化させながらリサンプリングを行なう。この場合、各仮説の入力訓練データに対する予測誤差の関数として各仮説の重みを計算し、それらの重みによる重み付き平均をもって最

$$M(x) = \sum_{H_i(x)=Y_{\max}} w_i - \max_{Y \neq Y_{\max}} \sum_{H_i(x)=Y} w_i \quad \dots \dots (2)$$

ここで、 $H_i(x)$ は仮説 H_i のデータ点 x に対する予測値、 w_i は、仮説 H_i に付された重みを指す。また、 y_{\max}

$$y_{\max} = \arg \max_Y \sum_{H_i(x)=Y} w_i \quad (3)$$

となる。この場合、マージン最小の点が情報量最大の点と推定される。なお、各仮説の重みを導入しない場合においても、マージンの概念は各仮説の重みを1とすることによって拡張可能なので、上記の手順によってマージン最小の点を情報量最大の点として選択することが可能である。

$$V(x) = \sum_{i=1, \dots, T} w_i |H_i(x) - \overline{H}_i(x)| \quad (4)$$

ここで、 $H_i(x)$ にバーを付したものは x のラベル予測値の平均値である。分散値の重み付き平均値の大きいが候補データの情報量は大きいとみなす。

【0024】以上、各仮説の予測値の分散値やマージン値によって情報量を算出することを説明したが、各データ点に対する正しいラベル値を用いて算出可能な、仮説

$$M'(x) = \sum_{H_i(x)=Y^*} w_i - \max_{Y \neq Y^*} \sum_{H_i(x)=Y} w_i \quad (5)$$

ここで、 y^* はデータ x の真のラベル値を表す。

【0025】また、発見すべき知識表現が実数関数の場合には、各仮説のデータ点 x に対する予測誤差の重み付き平均値 $V'(x)$ として算出することが可能である。

$$V'(x) = \sum_{i=1, \dots, T} w_i |H_i(x) - y^*| \quad (6)$$

以上の定義によれば、 $M'(x)$ は $M(x)$ の定義中の予測値モードを真のラベル値で置き換えることにより得られ、 $V'(x)$ は $V(x)$ の定義中の予測平均値を真の値で置き換えることにより得られる。一般には、二乗誤差以外の誤差の測度、例えば絶対誤差を用いることも可能である。

終仮説による予測を行なう。そして、データ選択の際の情報量の推定にもこれらの重みを用いる。例えば、発見すべき知識表現が多値関数の場合には、各候補データに対して、各予測値を予測する仮説の重みの総和を算出し、もっとも大きい重みの総和と次に大きい重みの総和との差分（以下、この量をマージンと呼ぶ）を用いて、情報量を測ることができる。

【0022】上記の手順をより詳細に説明する。データ点 x に対するマージン $M(x)$ は以下のように定義される。

【数1】

x は重み総和最大の予測値とする。即ち、

【数2】

【0023】発見すべき知識表現が実数関数の場合には、各候補データに対して、予測値の分散値の重み付き平均値を用いて、情報量を測ることができる。分散値の重み付き平均 $V(x)$ は、次式で表される。

【数3】

30 の予測誤差を利用することも可能である。例えば、発見すべき知識表現が多値関数の場合、(2) 式のマージンの定義を修正して、次式で表される誤差マージン M'

(x) を定義して誤差マージン最小のデータ点を情報量最大とすることができる。

【数4】

例えば、誤差として二乗誤差を用いれば $V'(x)$ は 40 次式から算出される。

【数5】

【0026】さらに、上記の上位アルゴリズムにおいて、メインメモリ量に限度があり、繰り返し選択したデータの蓄積がその限度を超過してしまう場合や、繰り返し選択し、メインメモリに格納したデータの一部が学習に不要となる場合、あるいはそれらデータを学習するに 50 は計算時間がかかり全てを学習することが実時間で不可

能な場合に、上記誤差マージンや $V'(x)$ を利用して、メインメモリのデータを削除することが可能である。これは、メインメモリに格納されたデータのうち、情報量の小さなデータを削除することにより達成される。例えば、発見すべき知識表現が多値関数の場合には、上記誤差マージンを利用して、誤差マージンが比較的大きな点を1つまたは複数、メインメモリから削除する。また、発見すべき知識表現が実数関数の場合には、上記 $V'(x)$ を算出し、 $V'(x)$ の小さなデータから削除する。

【0027】以下、図面を参照してこの発明の実施の形態について説明する。図1は、同実施形態による知識発見方式のプログラム100の構成を示す図である。プログラム100は、下位学習アルゴリズムと訓練データとして学習させるデータを入力する入力部107と、精度増強部103と、入力部107へ入力されたラベル未知のデータに対しラベル値を予測する予測部102と、予測部102が予測したラベル値を出力する出力部101と、精度増強部103における繰り返し処理を制御する制御部106とから構成される。

【0028】図1の各部について説明する。精度増強部103は、高速メインメモリ120に格納されているデータから再サンプリングによって作成された複数の部分サンプルを下位学習アルゴリズムに学習させ、複数の仮説を得る学習手段104と、学習手段104により得られた複数の仮説を用いて、大容量データ格納装置より読み出された候補データ各々の情報量の推定を行ない、情報量の大きいデータのみを高速メインメモリ120に追加格納するデータ選択手段105からなる。大容量データ格納装置130、高速メインメモリ120およびCPU（中央演算処理装置）110は、プログラム100を実行する計算機を構成する要素である。また、精度増強部103のその他の動作例を説明する。精度増強部103は、メインメモリ120に格納されているデータから学習を行い、学習した仮説をメインメモリ120もしくは大容量データ格納装置130に格納する学習手段104と、大容量データ格納装置130より読み出された候補データに対し、過去に学習された仮説を用いて情報量の推定を行い、情報量の大きなデータのみをメインメモリ120に格納するデータ選択手段105からなる。

【0029】次に、上記構成による知識発見方式の動作を図面を参照して説明する。図2は、本発明の実施例の動作の流れを示す図である。先ず、下位学習アルゴリズム、例えば、決定木を学習するアルゴリズムなどを入力部107へ入力する（ステップS201）。次に、ステップS202で、学習手段104は、この時点で高速メインメモリ120に格納されているデータからリサンプリングによって部分データ集合を作成する。この部分データ集合を入力部107を介して入力し、前記下位学習アルゴリズムにより学習させ、仮説を得る（ステップS

203）。

【0030】次に、リサンプリング回数 i と停止条件となるリサンプリング回数 T を比較し、リサンプリング回数 i が所定の回数 T を越えないとき（ステップS204；NO）、ステップS202に戻り、上記の処理を繰り返す。リサンプリング回数 i が所定の回数 T を越えると（ステップS204；YES）、データ選択手段105は、最終的に得られた前記仮説を用いて、大容量データ格納装置130から読み出された候補データの情報量の推定を行い情報量の大きいデータを選択する（ステップS205）。データ選択手段105により選択されたデータを高速メインメモリ120に格納し、既に格納してあるデータに加える（ステップS206）。

【0031】次に、ステップS207において、ステージ数 j を停止条件となるステージ数 N と比較し、ステージ数 j が所定の回数 N を越えないとき（ステップS207；NO）、ステップS202から処理を繰り返す。ステージ数 j が所定の回数 N を越えたとき（ステップS207；YES）、知識発見の過程を終了し、得られた規則を最終仮説として出力する。図2と同様に、図3も本発明の実施例の動作の流れを示す図である。図3においては、まず下位学習アルゴリズムを入力部107へ入力する（ステップS301）。次に、現在メインメモリに格納されている学習データを入力して前記下位学習アルゴリズムに学習させ、学習した仮説をメインメモリ120もしくは大容量データ格納装置130に格納する（ステップS302）。次に、データ選択手段105は、これまで過去に得られた仮説を用いて、大容量データ格納装置130から読み出された候補データの情報量を推定し、情報量の大きなデータを選択し（ステップS303）、メインメモリに格納する（ステップS304）。最後に、ステップS305において、ステージ数 j を停止条件となるステージ数 N と比較し、ステージ数 j が所定の回数 N を越えない時（ステップS305；NO）、ステップS302から処理を繰り返す。ステージ数 j が所定の回数 N を越えた時（ステップS305；YES）、知識発見の過程を終了し、得られた規則を最終仮説として出力する。

【0032】

40 【発明の効果】以上説明したように、本発明の知識発見方式によれば、ディスクなどの大容量データ格納装置に格納された膨大な量のデータから情報量の大きい部分だけを選択的にサンプリングして得られた比較的少数のデータから効率的に、高精度で未知データのラベル予測を行なう規則を発見することが可能になり、データマイニングの実効性を高めるという効果が得られる。また、データの選択的サンプリングにより少数のデータをメインメモリに読み込み知識発見を行うことからメインメモリに制限のある計算機でもデータマイニングを行うことができるという効果が得られる。

【図面の簡単な説明】

【図1】 この発明の一実施の形態の構成を示す図である。

【図2】 同実施の形態の動作の流れを示す図である。

【図3】 同実施の形態の動作の流れを示す図である。

【符号の説明】

100 知識発見プログラム

101 出力部

102 予測部

103 精度増強部

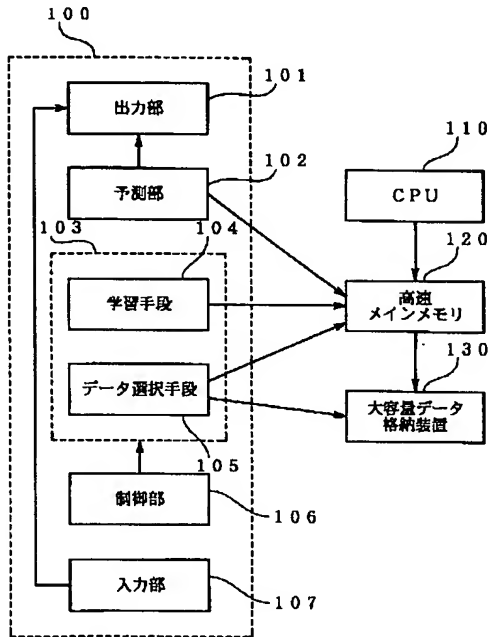
104 学習手段

105 データ選択手段

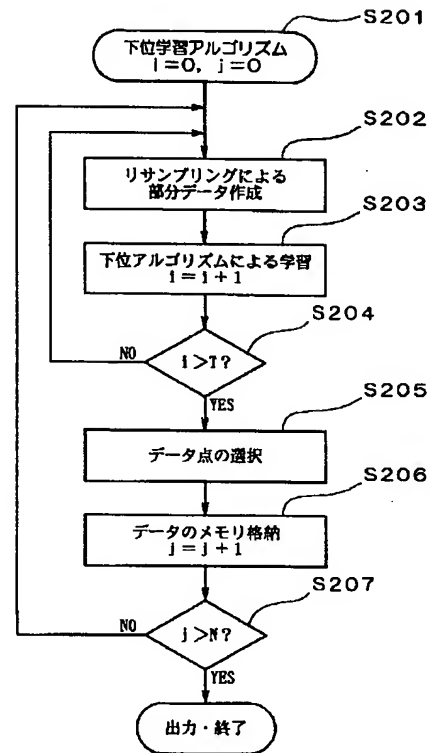
106 制御部

107 入力部

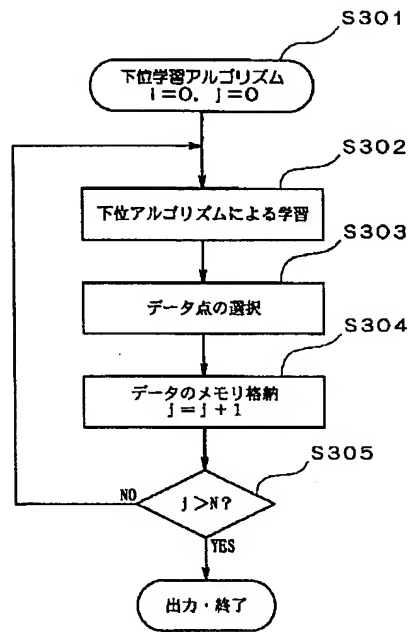
【図1】



【図2】



【図 3】



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☐ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☒ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.